

(19)



JAPANESE PATENT OFFICE

PATENT ABSTRACTS OF JAPAN

(11) Publication number: 10124533 A
(43) Date of publication of application: 15.05.1998

(51) Int. Cl. G06F 17/30
G06F 12/00

(21) Application number: 09121367
(22) Date of filing: 13.05.1997
(30) Priority: 13.05.1996 US 96 644599

(71) Applicant: LUCENT TECHNOLOG INC
(72) Inventor: GANGULY SUMIT
GIBBONS PHILLIP B
MATIAS YOSHI
SILBERSCHATZ ABRAHAM

(54) METHOD FOR EVALUATING BIAS
PREVENTION CONNECTION SIZE

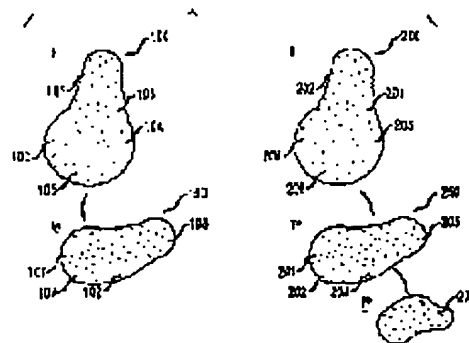
(57) Abstract:

PROBLEM TO BE SOLVED: To obtain a method for evaluating a data base question sensitive to the possibility of the bias of data by applying an evaluating procedure to each defined group, and deciding the whole evaluation of a data base.

SOLUTION: A first data base R and a second data base T are respectively provided with plural data items 101-103 and 201-203 having different coupling attribute values. The evaluation is operated to the data items 101-103 and 201-203 having the specific coupling attribute densely occupying the both data bases R and T. Then, the evaluation for suppressing the influence of the dense data items 101-103 having the specific coupling attribute value in the first data base R

is operated, and finally the evaluation for suppressing the influence of the dense data items 201-203 having the same coupling attribute value in the second data base T is operated. Then, the dense - dense and non-dense - any relation is defined by using the threshold value of each data base R and T.

COPYRIGHT: (C)1998.JPO



(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号

特開平10-124533

(43) 公開日 平成10年(1998) 5月15日

(51) IntCl ⁶	識別記号	F I
G 0 6 F 17/30		G 0 6 F 15/403
12/00	5 1 3	12/00
		15/403
		3 3 0 Z
		5 1 3 Z
		3 4 0 D

審査請求 未請求 請求項の数13 O L (全 8 頁)

(21) 出願番号 特願平9-121367
(22) 出願日 平成9年(1997) 5月13日
(31) 優先権主張番号 08/644599
(32) 優先日 1996年5月13日
(33) 優先権主張国 米国 (US)

(71) 出願人 596092698
ルーセント テクノロジーズ インコーポ
レーテッド
アメリカ合衆国, 07974-0636 ニュージ
ャーシー, マレイ ヒル, マウンテン ア
ヴェニュー 600
(72) 発明者 サミット ガングリー
インド国 462024 ボーバル, シャクティ
ナガー, セクター-2 192
(74) 代理人 弁理士 岡部 正夫 (外9名)

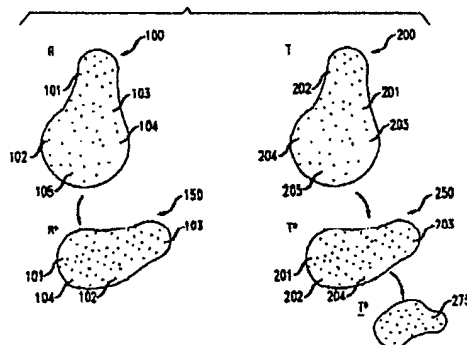
最終頁に続く

(54) 【発明の名称】 偏り防止結合サイズ評価方法

(57) 【要約】

【課題】 2つのデータベースT及びRの質問サイズの
評価方法を提供する。

【解決手段】 この方法は、データベースを稠密または
疎であるとして類別するためにスレシヨールドを使用す
る。そこで、2つのデータベースに稠密-稠密手順が適
用され、稠密-稠密評価 (A_{oo}) を作り出す。データベ
ースTからくるデータアイテムを抑制する疎-何れか手
順が行なわれ、第1の疎-何れか評価 (A_{o1}) を作り出
す。次いで、データベースRから稠密なデータアイテム
を抑制することによって、第2の疎-何れか評価
(A_{oo}) が作り出される。最後に、稠密-稠密評価、第
1の疎-何れか評価及び第2の疎-何れか評価を結合す
ることにより、質問サイズ評価が作り出される。



【特許請求の範囲】

【請求項1】 少なくとも2つのデータベースR及びTの等結合を評価することによりデータを管理する方法であって、

Rの稠密なデータアイテムとTの稠密なデータアイテムの質問サイズを評価することにより、稠密-稠密評価を作り出す工程と、

Rの稠密なデータアイテムを抑制する質問サイズを評価することにより、第1の疎-何れか評価を作り出す工程と、

Tの稠密なデータアイテムを抑制する質問サイズを評価することにより、第2の疎-何れか評価を作り出す工程と、

前記稠密-稠密評価と、前記第1の疎-何れか評価と、前記第2の疎-何れか評価を結合することにより、データベースR及びTの等結合のサイズの評価を作り出す工程とからなる方法。

【請求項2】 請求項1記載の方法において、前記結合工程は、前記稠密-稠密評価と、前記第1の疎-何れか評価と、前記第2の疎-何れか評価を加算することにより行われる方法。

【請求項3】 請求項1記載の方法において、前記結合工程は、前記稠密-稠密評価と、前記第1の疎-何れか評価と、前記第2の疎-何れか評価を平均することにより行われる方法。

【請求項4】 請求項1記載の方法において、前記結合工程は、前記稠密-稠密評価と、前記第1の疎-何れか評価と、前記第2の疎-何れか評価のうちの1つを最大値として選択することにより行われる方法。

【請求項5】 少なくとも2つのデータベースR及びTのデータベース質問サイズを評価することによりデータを管理する方法であって、

データベースRにおけるある属性値を有する稠密なデータアイテムと前記属性値を有するTの稠密なデータアイテムとの評価を行なうことにより、稠密-稠密評価を作り出す工程と、

前記データベースRの稠密なデータアイテムを抑制する評価を行なうことにより、Rにおける疎-何れか評価を作り出す工程と、

Tの稠密なデータアイテムを抑制する評価を行なうことにより、Tにおける疎-何れか評価を作り出す工程と、前記稠密-稠密評価と、前記Tにおける疎-何れか評価と、前記Rにおける疎-何れか評価を結合することにより、前記データベースR及び前記データベースTのデータベース質問サイズを評価する工程とからなる方法。

【請求項6】 請求項5記載の方法において、前記稠密-稠密評価は、

前記データベースRからデータアイテムをサンプリングすることにより、サンプルR*を作り出し、

前記データベースTからデータアイテムをサンプリング

することにより、サンプルT*を作り出し、

前記サンプルR*におけるある結合属性値(v)を有する多数のデータアイテムを決定することにより、Rにおいて多数の前記結合属性値(v)を作り出し、

前記サンプルT*における前記結合属性値(v)を有する多数のデータアイテムを決定することにより、Tにおいて多数の前記結合属性値(v)を作り出し、

前記結合属性値(v)の各々に関して、前記結合属性値(v)のサブ結合のサイズの中間の稠密-稠密評価を決定し、

前記結合属性値(v)の各々に関して中間の稠密-稠密評価を加算し、

前記サブ結合属性値(v)の前記サイズの中間の稠密-稠密評価を見積もることにより行なわれる方法。

【請求項7】 請求項6記載の方法において、前記中間の稠密-稠密評価は、スレシヨールド値を決定し、Tにおける多数の前記結合属性値(v)及びRにおける多数の前記結合属性値(v)を前記スレシヨールド値と比較した後に行われる方法。

【請求項8】 請求項7記載の方法において、前記中間の稠密-稠密評価は、T*における多数の前記結合属性値(v)及びR*における多数の前記結合属性値(v)が共に前記スレシヨールド以上であることを決定した後に行われる方法。

【請求項9】 請求項7記載の方法において、前記中間の稠密-稠密評価は、T*における多数の前記結合属性値(v)及びR*における多数の前記結合属性値(v)が共に前記スレシヨールドに等しいことを決定した後に行われる方法。

【請求項10】 請求項5記載の方法において、前記第1の疎-何れか評価は、

前記データベースTからデータアイテムをサンプリングすることにより、サンプルT*を作り出し、

Rにおけるその数がスレシヨールド以上であるT*のデータアイテムを抑制し、

前記結合属性値(v)の各々について、Tにおいて疎である、前記結合属性値(v)を有するデータアイテムにより中間の疎-何れか評価を計算し、

前記結合属性値(v)の各々について中間の疎-何れか評価を加算し、

前記中間の疎-何れか評価を見積もることにより、RにおけるTの疎-何れか評価を作り出すことにより行なわれる方法。

【請求項11】 請求項10記載の方法において、T*のデータアイテムを抑制する前記工程は、稠密な結合属性値(v)を有する、Rにおける多数のデータアイテムを決定することにより行なわれる方法。

【請求項12】 請求項10記載の方法において、T*のデータアイテムを抑制する前記工程は、RからランダムサンプルR*をとり、R*に表われる各結合属性値

(v) に関して、 T^* から結合属性値 (v) を有する全てのデータアイテムを削除することにより行なわれる方法。

【請求項13】 請求項5記載の方法において、前記第2の疎一何れか評価は、前記データベースRからデータアイテムをサンプリングすることにより、サンプル R^* を作り出し、前記結合属性値 (v) を有する、Rにおける多数のデータアイテムを計算するコストを決定し、 R^* の各データアイテムについて、前記結合属性値 (v) を有するTのデータアイテムの数を決定し、前記結合属性値 (v) の各々について、Rにおいて疎である、前記結合属性値 (v) を有するデータアイテムにより中間の疎一何れか評価を計算し、前記結合属性値 (v) の各々について中間の疎一何れか評価を加算し、前記中間の疎一何れか評価を見積もることにより、TにおけるRの疎一何れか評価を作り出すことによって行なわれる方法。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】 本発明は、データベース質問評価に関する。

【0002】

【従来の技術及び発明が解決しようとする課題】 コンピュータのまん延につれて、コンピュータデータベースも増加した。近頃のデータベースのサイズは非常に大容量になることがあり、データベースの中には数百乃至数十億のデータアイテムを保持しているものがある。これらのデータベースの内の1つのデータベースの質問では、あらゆるデータアイテムが質問に合う可能性があり、あらゆるデータアイテムを比較しなければならないことがある。したがって、これらのデータベースのサイズが増すにつれて、質問を実行するコストも増加している。

【0003】 データベースは1つ以上のテーブルからなり、各テーブルは数百乃至数十億のデータ記録を保持している。各データ記録は、情報が入っている1つ以上のフィールドを含む。これらのフィールド内の情報に基づいて、記録をいくつかのタイプのうちの1つのタイプとして類別することができる。例えば、テーブルは人間の記録を入れることができ、各記録は、人名を与えるフィールドと、好きなスポーツを与えるフィールドを有する。各記録は好きなスポーツで類別することができるので、タイプ“野球”、“フットボール”または“ホッケー”とすることができる。

【0004】 データベースのテーブルの質問に早く応答することができるのが望ましい。共通の質問の1つは等結合質問である。2つのテーブルR及びTの等結合質問の結果は、1つの記録がRからかつもう1つの記録がTからのものであつていずれも同タイプのものである記録ペアの全てからなるテーブルとなる。例えば、テーブル

Rが男性及び彼らが好きなスポーツの記録を含み、かつテーブルTが女性及び彼女らが好きなスポーツの記録を含む場合は、テーブルR及びTの等結合は、好きなスポーツが同じ男性と女性の記録ペアの全てを含むテーブルとなる。

【0005】 等結合の結果を計算するのは、コストがかかり過ぎることになる可能性がある。例えば、一方のテーブルがn個の記録を持ち、かつ他方のテーブルがn個の記録を持っている場合は、結果の計算は n^2 の記録ペアの比較を必要とし得る。 n^2 の比較をそれぞれ実行すると、質問コストが増加する。したがって、等結合質問のコストを下げるのが好適である。

【0006】 大容量データベースにおける費用のかかる質問の実行のコストを減らす必要性の結果として、データベース評価の分野がポピュラーになってきた。データベース評価において、評価は、データベースにおける質問の可能な出力（質問の評価と呼ばれる）で作られる。したがって、質問の評価は、質問を行なう前に計算される。その結果、質問を続けてコストを負うべきかまたはこの特定の質問を取り消すべきかに関して決定することができる。データベース評価の問題点は、評価が正確かつ計算が能率的になるように、特定の質問の評価を計算することを伴うことである。

【0007】 このデータベース評価の問題点を解決する試みが以前に行われた。パラメータ法と呼ばれる手法は、データを取り、データベースにおいてこのデータと既知のデータモデルを比較するのを試みる。データモデルの作用と性質がデータベースのデータと同じであることが原理になっている。しかしながら、パラメータ法では、データが既知のデータモデルにどのくらい近似しているかについての仮定が行われる必要があり、近似（データの適合）は結果の正確さをはなはだしく変えてしまうことがある。他の種類の手法は、データベースにおけるデータアイテムのサンプル（小さな組）を取り、これらのサンプルに基づいて評価を行なう。この種の手法はサンプリング法として知られている。サンプルは、通常、個々のデータベースの中からとられ、次いで合成されて質問評価を作り出す。サンプリング法は、他のデータベース評価法に勝る利点があることが証明された。パラメータ法と違って、サンプリング法は、データの適合について仮定がより少なくな行われるべきことを要する。さらに、サンプリング法は、常に、統計的確実性と呼ばれる正確さの確実性を有する。たいいていのサンプリング法の統計的確実性は、典型的に、90%乃至99%の範囲になるのを目指している。

【0008】 サンプリング法もパラメータ法もデータベース評価には好適ではない。しかしながら、サンプリング法は、データベースのデータアイテムが偏っている場合、パラメータ法より正確な結果を提供することができる。例えば、上記に説明したデータベースが同型でない

(同型でないデータベースとは、データベース内のデータアイテムが異なるタイプからなる同等でない混合になっているデータベースである) 場合、パラメータ法は良好な性能を示すことができない。他のタイプより1つのタイプのデータアイテムがかなり多い(例えば20%以上) 場合、データベースは偏っているといわれる。偏ったデータベースは、パラメータ法やサンプリング法を実行する場合に問題があることがわかっている。

【0009】 サンプルが多く偏ったデータアイテムを含んでいるかまたはあまり偏っていないデータアイテムを含んでいるかに依存して、偏ったデータの影響は、サンプリング法の結果に劇的に影響を与えることがある。サンプリング法を使用する場合、偏ったデータの影響を考慮することができない大量のサンプルがテーブルT及びテーブルRからとられる。したがって、偏ったデータを考慮してサンプリング法を実行するのは不具合がある。

【0010】

【課題を解決するための手段】 本発明は、データの偏りの可能性に敏感な、データベース質問の評価方法を実行する。データベース質問の総合評価は、3つの個別評価を合わせて行われる。まず、評価は、稠密に両データベースR及びTを占める、特定の結合属性を有するデータアイテムについて行われる。次に、第1のデータベースRにある、特定の結合属性値を有する稠密なデータアイテムの影響を抑制する評価が行われ、最後に、第2のデータベースTにある、同じ結合属性値を有する稠密なデータアイテムの影響を抑制する評価が行われる。

【0011】 本発明は、データベース質問の評価へ広範囲にわたるアプローチをとっている。この広範囲にわたるアプローチは、2つの異なるデータベースR及びTのデータアイテムの関係を確立することによって作り出される。この関係は、サブ結合の収集からなる2つの部分に分かれたグラフとして知られている。各サブ結合は、この特定の結合属性値を有するデータアイテムの全ペアからなる。

【0012】 各データベースのランダムなサンプルが収集される(例えば、R* 及びT* は、それぞれデータベースR及びTのランダムなサンプルを示す)。次いで、ランダムサンプルR* 及びT* におけるサブ結合は、特定の結合属性値のデータアイテム数がスレシヨールド値以上か以下かに基づいて、特定の結合属性値のデータアイテムの稠密な母集団または特定の結合属性値のデータアイテムの疎な母集団を持つように評価される。本発明の方法では、データアイテムは、その結合属性値が稠密ならば稠密になり、その結合属性を有するデータアイテムの母集団が疎ならば疎になることがわかる。スレシヨールド値は、データベースのデータアイテム数の平方根として定義される。サブ結合の各ペアのデータアイテムが共に稠密な場合は、サブ結合は稠密-稠密と呼ばれ

る。サブ結合の各ペアにおいて、一方のデータアイテムが疎なものとして類別され、他方のデータアイテムが稠密または疎なものとして類別されている場合は、サブ結合は、疎-何れか(例えば、稠密-疎、疎-稠密または疎-疎) として類別される。

【0013】 次いで、一連の3つの評価がデータベースで行われる。手順はまず、稠密-稠密評価(A_d) を決定するためにランダムサンプルに適用される。次いで、Rの結合属性値を有する稠密データアイテムを抑制する疎-何れか(A_{S1}) 評価が行われ、次に、Tの結合属性値を有する稠密データアイテムを抑制する疎-何れか(A_{S2}) 評価が行われる。最後に、稠密-稠密評価と、Rの稠密データアイテムを抑制する疎-何れか評価と、Tの稠密データアイテムを抑制する疎-何れか評価とを結合することによって、データベース質問評価(A) が作られる。

【0014】 稠密-稠密評価と疎-何れか評価は、特定のデータベース質問に対して、他のランダムサンプル(例えばT*) と適合する可能性のある一方のランダムサンプル(例えばR*) における結合属性値を有するデータアイテムを評価することにより行われる。上述のように、特定の結合属性値を有するデータアイテムの適合はこの方法ではサブ結合と呼ばれる。したがって、各結合属性値に関して、その値に関連したサブ結合のサイズが評価される。次いで、サブ結合の評価の和が全ての結合属性値に関して加算され、これらの和(例えば、稠密-稠密評価、疎-何れか評価) は組み合わせられてデータベース質問評価を作り出す。

【0015】 本発明の目的、利点及び新規な特徴は、添付図面に関して読まれる以下の詳細な説明からより十分に明らかになるだろう。

【0016】

【発明の実施の形態】 本発明では、等結合データベース質問サイズの評価方法が実行される。まず、質問基準が定義される。各データベースからのデータアイテムは質問基準にしたがってグループ分けされる。次いで、“稠密-稠密” 手順及び“疎-何れか” として知られる評価手順が、定義されたグループの各々に適用され、データベースの全体評価を決定する。

【0017】 詳細には、この方法は、2つのデータベースの稠密-稠密評価(A_d) を行なうことによってデータベースR及びTの等結合評価(A) を行なう。次いで、Rの稠密なデータアイテムを抑制する疎-何れか手順が行われ、第1の疎-何れか評価(A_{S1}) を作り出す。次いで、Tの稠密なデータアイテムを抑制する疎-何れか手順が行われ、第2の疎-何れか評価(A_{S2}) を作り出す。最後に、 A_d 、 A_{S1} 及び A_{S2} を組み合わせる(加算する) ことによって、2つのデータベースの等結合評価(A) が計算される。

【0018】 本発明の方法を示す手段として、図1は、

説明を容易にするために用いることができる概念的モデルを表わしている。図1において、100で示された第1のデータベース(R)及び200で示された第2のデータベース(T)が示されている。第1のデータベースR及び第2のデータベースTは共に、異なる結合属性値(例えば、それぞれ101, 102, 103及び201, 202, 203)を有する多数のデータアイテムを有している。また、2つのデータベースサンプルが示されている。第1のランダムサンプル(R*)は150で示され、第2のランダムサンプル(T*)は250で示されている。両ランダムサンプル150及び250は、原データベースR及びTからのデータアイテムのランダムサンプリングを含んでいるだろう。したがって、160で示された第1のランダムサンプル(R*)はデータアイテム101, 102及び103を含み、250で示された第2のランダムサンプル(T*)はデータアイテム201, 202及び203を含む。

【0019】図1に表わされた概念的モデルの直観的な理解を進展させるために、データベースRのデータアイテムは、好きなスポーツを持っている男性の母集団を表わすと仮定する(好きなスポーツは結合属性値となるだろう)。したがって、データアイテム101, 102及び103は各々、特定のスポーツが好きな特定の男性を表わす。次に、データベースTは、好きなスポーツを持っている女性の母集団を表わすと仮定する。したがって、各データアイテム201, 202及び203は、特定のスポーツが好きな女性の男性を表わす。さらに、簡単にするために、安静または女性の各々が好きなスポーツ(結合属性値)は、ベースボール、フットボールまたはホッケーのどれかであると仮定する。したがって、データベースRは、好きなスポーツがベースボール、フットボールまたはホッケーである男性の母集団を表わし、データベースTは、好きなスポーツがベースボール、フットボールまたはホッケーである女性の母集団を表わす。

【0020】データベースR及びTの可能な質問(等結合)は、(1)ベースボールが好きな男性と女性、(2)フットボールが好きな男性と女性、(3)ホッケーが好きな男性と女性に適合することができる。このタイプの質問は、評価技術なしに行われる場合は3ステップで達成することができる。まず、データベースRは、好きなスポーツがベースボールである全男性についてサーチされるだろう。次に、データベースTは、好きなスポーツがベースボールである全女性についてサーチされるだろう。最後に、好きなスポーツがベースボール(フットボール及びホッケー)である男性と女性の好一対が作られる。データベースR及びTが共にnデータアイテムを含んでいる場合、好きなスポーツが同じである男生と女性を組み合わせるのは n^2 の作業になる。

【0021】本発明の方法論では、ランダムサンプルR

*及びT*は、それぞれデータベース(R)及び(T)から取られる。次いで、サンプルは結合属性値で類別される。結合属性値でデータアイテムを類別するために、好きなスポーツは各々、それと関連する番号(例えば、ベースボール“1”、フットボール“2”及びホッケー“3”)を持っていると仮定する。そこで、特定のスポーツが好きな男性と女性の好一対は、図1の2つの部分に分かれたグラフで表わすことができる。図1Bにおいて、図1AのランダムサンプルR*及びT*からの各データアイテムは、その上に好きなスポーツの番号を伴って示されている。例えば、データアイテム101は、(データアイテム101上の数字が1なので)ベースボールが好きな、ランダムサンプルR*からの男性である。図1Bのデータアイテム201はランダムサンプルT*からのデータアイテムである。したがって、図1Bのデータアイテム201は、(その上に2があるので)フットボールが好きな女性である。

【0022】そこで、特定のスポーツが好きな男性と女性をペアにして、2つに別れたグラフとして表わすことができ、このグラフには、各男性及び女性を表わすノードと、男性と女性が好きなスポーツが同じ場合の、男性を表わすノードと女性を表わすノード間の線(データアイテム間に引かれた線)とがある。図1Bにおいて、データアイテム101, 102及び103は、(それぞれ、1, 3及び2で表わされるスポーツが好きな)3人の男性の記録を表わし、データアイテム201, 202及び203は、(それぞれ、2, 1及び2で表わされるスポーツが好きな)3人の女性の記録を表わす。302, 304, 306で示された線は、103で示された男性(ベースボールが好きな男性)を(ベースボールが好きな)3人の女性とペアにする。

【0023】一般に、各結合属性値(例えば、ベースボール、フットボール、ホッケー)について、この属性値(例えば1, 2, 3)を有する男性の各々を同じ属性値(例えば1, 2, 3)を有する女性とペアにする線がある。ある特定の結合属性値のデータアイテムと線を1組にして、サブ結合として示されている。正確には、各結合属性値(例えば1で示されたベースボール、2で示されたフットボール及び3で示されたホッケー)に対して1つのサブ結合がある。例えば、ベースボールが好きな男性と女性のサブ結合は、データアイテム103, 104と、線302, 303, 304, 305, 306, 307と、データアイテム201, 203及び204とからなるだろう。

【0024】個々のデータベースのスレシールド値を用いて、本発明は、稠密-稠密及び疎-何れか(すなわち、稠密-疎、疎-稠密及び疎-疎)のような関係を定義する。スレシールド値は、ランダムサンプルを稠密なデータタイプの母集団の表現にすることが可能などんな値でも良いことがわかる。この関係は次のように定義

される。

【0025】稠密-稠密データベース(R)におけるある結合属性値を有するデータアイテムの数とデータベース(T)における同じ結合属性値を有するデータアイテムの数が各々、平方根(n)より大きいまたは等しい関係。ここで、nはデータベースのデータアイテム数である。

【0026】稠密-疎データベース(R)におけるある結合属性値を有するデータアイテムの数が平方根(n)より大きいまたは等しく、かつデータベース(T)におけるデータアイテムの数が平方根(n)より小さい関係。ここで、nは各データベースのデータアイテム数である。

【0027】疎-稠密データベース(R)における特定の結合属性値を有するデータアイテムの数が平方根(n)より小さく、かつデータベース(T)におけるデータアイテムの数が平方根(n)より大きいまたは等しい関係。ここで、nは各データベースのデータアイテム数である。

【0028】疎-疎データベース(R)及びデータベース(T)における特定の結合属性値を有するデータアイテムの数が平方根(n)より小さい関係。ここで、nは各データベースのデータアイテム数である。

【0029】本発明の方法論では、ある結合属性値vを有する1組のデータアイテムについて、“mult_T(v)”は、その結合属性値を有するTにおけるデータアイテムの数と定義される。Rにおいて、“mult_R(v)”がnの平方根より大きいまたは等しければ、結合属性値は稠密であると定義され、また、“mult_R(v)”がnの平方根より小さければ、結合属性値が疎であると定義される。さらに、サブ結合(v)は、vがR及びTの両方で稠密ならば、稠密-稠密サブ結合となる。サブ結合(v)は、vがR及びTの両方で疎ならば、疎-疎サブ結合となる。

【0030】本発明の方法論では、nは各関係におけるデータアイテム数であり、m₁は稠密-稠密手順のサンプルサイズであり、M₂は疎-何れか手順のサンプルサイズであり、δは稠密-稠密手順で使用されるスレシールド値である。上記に定義された変数を仮定すれば、各評価は次のように定義することができる。

【0031】 $A_d := f_d(n, m_1, \delta)$

$A_{s1} := f_s(R, T, n, m_1)$

$A_{s2} := f_s(T, R, n, m_1)$

$A := A_d + A_{s1} + A_{s2}$

【0032】ここで、 $f_d(x)$ 及び $f_s(x)$ はxの関数であり、Aは総合データベース評価であり、 A_d は稠密-稠密データベース評価であり、 A_{s1} は、Rの稠密なデータアイテムを抑制する疎-何れかデータベース評価であり、 A_{s2} はTの稠密なデータアイテムを抑制する疎-何れか評価である。 $A < n \log n$ ならば、この方

法は、質問サイズの上限である健全な限界(S)： $= n \log n$ も提供する。開示された実施例では、 $m_1 = (\sqrt{n} + \log n) * \log n$ 、 $m_2 = \sqrt{n} + \log n$ 、 $\delta = \log n$ となる。

【0033】本発明では、各データベースのランダムサンプルがとられる。2つのデータベースからとられたランダムサンプルはほとんど $m_1 + m_2$ であり、 m_1 は、稠密-稠密評価を行なう場合に(個別的に)R及びTからサンプリングされたデータアイテム数であり、 m_2 は、疎-何れか評価を行なう場合に(個別的に)R及びTからサンプリングされたデータアイテム数である。サンプルがとられると、稠密-稠密評価(A_d)が引き出され、Rの稠密なデータアイテムを抑制する疎-何れか評価(A_{s1})が計算され、Tの稠密なデータアイテムを抑制する疎-何れか評価(A_{s2})が計算される。最後に、 A_d 、 A_{s1} 及び A_{s2} を結合することにより、質問サイズの評価を行なうことができる。

【0034】稠密-稠密

この方法論、稠密-稠密手順の第1のステップは、以下のステップを含む。

1. データベースR及びTからランダムサンプルR*及びT*がとられる。ランダムサンプルR*及びT*は各々サイズ m_1 からなる。
2. V^* はR*及びT*の両方におけるデータアイテムの1組の結合属性値とする。
3. 各値 $v \in V^*$ について、 $\text{mult}_{R^*}(v)$ を決定する。
4. 各値 $v \in V^*$ について、 $\text{mult}_{T^*}(v)$ を決定する。
5. 各値 $v \in V^*$ について、 $\text{mult}_{R^*}(v) \geq \delta$ かつ $\text{mult}_{T^*}(v) \geq \delta$ ならば、 $A' := A' + \text{mult}_{R^*}(v) * \text{mult}_{T^*}(v)$ となる。ここで、 A' は中間の稠密-稠密評価(A' は初期にゼロに設定される)であり、“*”は乗算を示すために使用される記号である。
6. 稠密-稠密評価 $A_d := (n/m_1)^2 * A'$ となる。

【0035】疎-何れか

疎-何れか評価を作り出す方法は以下の手順からなる。

1. データベースTからランダムサンプルT*がとられる。ランダムサンプルT*はサイズ m_2 からなる。
2. $\text{mult}_{T^*}(v)$ が稠密vについて多数の計算を必要とするならば、稠密vを抑制するために蓋然論的消去を使用する。
 - 2 a. データベースRからランダムサンプルR*がとられる。ランダムサンプルR*はサイズ m_2 からなる。
 - 2 b. R*に表われる各結合属性値について、T*から結合属性値vを有する全てのデータアイテムを削除する。
3. Rにおいて疎である残りの結合属性値に基づいて中

間の疎一何れか評価を計算する。

T* における各データアイテムyについて、

3 a. yと結合する、すなわちyと同じ結合属性値を有する、Rのデータアイテムの番号xを決定する。

3 b. $x < n/m_2$ ならば、 $A' := A' + x$ となる。ここで、A' は中間の疎一何れか評価である (A' は初期にゼロに設定される)。

4. $A_s := nA' / m_2$ となる。ここで、A_s は疎一何れか評価である。

【0036】本発明の方法論では、疎一何れか評価は2度行なわれる。疎一何れか手順が2回行なわれ、1回目に使用されなかったデータベースの稠密なデータアイテムが抑制される。例えば、データベースRの稠密なデータアイテムが抑制されたならば、1回目に疎一何れか手順が行なわれ、2回目にデータベースTの稠密なデータアイテムが抑制されるだろう。1回目に疎一何れか評価が行なわれると、その評価はA_{s1}で示され、2回目に (RとTの役割が反対にされて) 疎一何れか評価が行なわれると、その評価はA_{s2}で示される。稠密-稠密及び疎一何れかの評価の全てが計算されると、等結合評価は、 $A = A_s + A_{s1} + A_{s2}$ で示される。 $A < (n) \log(n)$ ならば、健全な限界 $S = (n) \log(n)$ の出力を計算することができる。ここで、Sは健全な限界である。すなわち、ハットは評価の統計的確信を維持する。さらに、等結合評価を計算するために、稠密-稠密評価、第1の疎一何れか評価及び第2の疎一何れか評価を追加するべく、等結合評価は、3つの評価を平均する

かまたは3つの評価の最大をとることにより近似することとできることがわかる。

【0037】さらに、本発明の方法論は、どんな数のデータベースにも適用することができる。例えば、3つのデータベースA、B及びCが含まれている場合は、この方法論は、まずデータベースA及びBに適用して中間の評価を作り出し、次いで、この方法論をCに用いて中間評価を結合することができる。

【0038】また、本発明は、データベース質問の正確で費用のかからない評価を作り出すために適用することができる。この方法は、質問最適化装置に、または並列もしくは分散データベースの多数のプロセッサにおける仕事量のバランスを取るのに必要なリソース割当ての決定時に有効である。さらに、より広範囲の規模に、開示されたこの方法を、会計監査や統計的研究等の大容量データベースアプリケーションにも適用することができる。

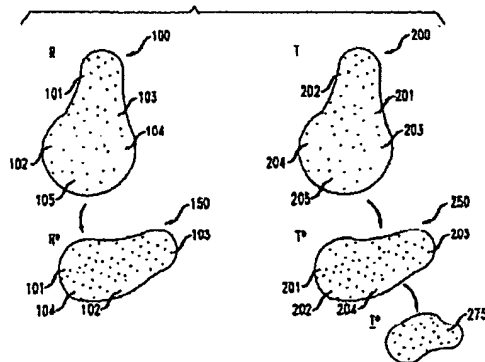
【0039】本発明のいくつかの実施例が開示されて説明されているが、本発明の精神または従属請求項の範囲から逸脱することなく種々の変更を行なうことができることがわかる。

【図面の簡単な説明】

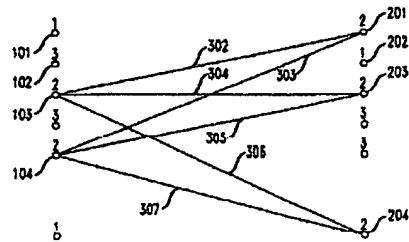
【図1A】母集団R及びTとランダムサンプルR* 及びT* の概念図を示す。

【図1B】ランダムサンプルR* 及びT* の選択されたデータアイテムと、結合属性値2のサブ結合の2つの部分に分かれたグラフを示す。

【図1A】



【図1B】



フロントページの続き

(72)発明者 フィリップ ビー. ギボンズ
アメリカ合衆国 07090 ニュージャージー
イ, ウェストフィールド, エンブリー コ
ート 201

(72)発明者 ヨッシ マティアス
アメリカ合衆国 20854 メリーランド,
ボタマック, ロザリング ドライヴ
11815

(72)発明者 アブラハム シルバーシャッツ
アメリカ合衆国 07901 ニュージャージー
イ, サミット, ニューイングランド アヴ
ェニュー 67エー

**This Page is Inserted by IFW Indexing and Scanning
Operations and is not part of the Official Record**

BEST AVAILABLE IMAGES

Defective images within this document are accurate representations of the original documents submitted by the applicant.

Defects in the images include but are not limited to the items checked:

- ☐ **BLACK BORDERS**
- ☐ **IMAGE CUT OFF AT TOP, BOTTOM OR SIDES**
- ☐ **FADED TEXT OR DRAWING**
- ☐ **BLURRED OR ILLEGIBLE TEXT OR DRAWING**
- ☐ **SKEWED/SLANTED IMAGES**
- ☐ **COLOR OR BLACK AND WHITE PHOTOGRAPHS**
- ☐ **GRAY SCALE DOCUMENTS**
- ☐ **LINES OR MARKS ON ORIGINAL DOCUMENT**
- ☐ **REFERENCE(S) OR EXHIBIT(S) SUBMITTED ARE POOR QUALITY**
- ☐ **OTHER:** _____

IMAGES ARE BEST AVAILABLE COPY.

As rescanning these documents will not correct the image problems checked, please do not report these problems to the IFW Image Problem Mailbox.